

# **TimeStitch: Exploiting Slack to Mitigate Decoherence in Quantum Circuits**

KAITLIN N. SMITH and GOKUL SUBRAMANIAN RAVI, University of Chicago, USA PRAKASH MURALI, Princeton University, USA JONATHAN M. BAKER, University of Chicago, USA

NATHAN EARNEST and ALI JAVADI-ABHARI, IBM Quantum, IBM T. J. Watson Research Center, USA

Center, USA

FREDERIC T. CHONG, University of Chicago, USA

Quantum systems have the potential to demonstrate significant computational advantage, but current quantum devices suffer from the rapid accumulation of error that prevents the storage of quantum information over extended periods. The unintentional coupling of qubits to their environment and each other adds significant noise to computation, and improved methods to combat decoherence are required to boost the performance of quantum algorithms on real machines. While many existing techniques for mitigating error rely on adding extra gates to the circuit [13, 20, 56], calibrating new gates [50], or extending a circuit's runtime [32], this article's primary contribution leverages the gates already present in a quantum program without extending circuit duration. We exploit circuit slack for single-qubit gates that occur in idle windows, scheduling the gates such that their timing can counteract some errors.

Spin-echo corrections that mitigate decoherence on idling qubits act as inspiration for this work. Theoretical models, however, fail to capture all sources of noise in Noisy Intermediate Scale Quantum devices, making practical solutions necessary that better minimize the impact of unpredictable errors in quantum machines. This article presents TimeStitch: a novel framework that pinpoints the optimum execution schedules for single-qubit gates within quantum circuits. TimeStitch, implemented as a compilation pass, leverages the reversible nature of quantum computation to boost the success of circuits on real quantum machines. Unlike past approaches that apply reversibility properties to improve quantum circuit execution [35], TimeStitch amplifies fidelity without violating critical path frontiers in either the slack tuning procedures or the final

This work is funded in part by EPiQC, an NSF Expedition in Computing, under grants CCF-1730082/1730449; in part by STAQ under grant NSF Phy-1818914; in part by NSF Grant No. 2110860; by the US Department of Energy Office of Advanced Scientific Computing Research, Accelerated Research for Quantum Computing Program; and in part by NSF OMA-2016136 and the Q-NEXT DOE NQI Center. This research used resources of the Oak Ridge Leadership Computing Facility at the Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725. KNS is supported by IBM as a Postdoctoral Scholar at the University of Chicago and the Chicago Quantum Exchange. GSR is supported as a Computing Innovation Fellow at the University of Chicago. This material is based upon work supported by the National Science Foundation under Grant # 2030859 to the Computing Research Association for the CIFellows Project. PM is supported by an IBM PhD Fellowship at Princeton University. FTC is Chief Scientist at Super.tech (a division of ColdQuanta) and an advisor to Quantum Circuits, Inc.

Authors' addresses: K. Smith, G. Ravi, J. Baker, and F. Chong, University of Chicago Department of Computer Science, 5730 S Ellis Ave, Chicago, IL 60637; emails: kns@uchicago.edu, gravi@uchicago.edu, jmbaker@uchicago.edu, chong@cs.uchicago.edu; P. Murali, Princeton University Department of Computer Science, 35 Olden St, Princeton, NJ 08544; email: prakashmurali@gmail.com; N. Earnest and A. Javadi, IBM Quantum, IBM Thomas J. Watson Research Center, 1101 Kitchawan Rd, Yorktown Heights, NY 10598; emails: nate@ibm.com, Ali.Javadi@ibm.com.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

 $\ensuremath{\text{@}}$  2022 Association for Computing Machinery.

2643-6817/2022/10-ART8 \$15.00

https://doi.org/10.1145/3548778

8:2 K. Smith et al.

rescheduled circuit. On average, compared to a state-of-the-art baseline, a practically constrained TimeStitch achieves a mean 38% relative improvement in success rates, with a maximum of 106%, while observing bounds on circuit depth. When unconstrained by depth criteria, TimeStitch produces a mean relative fidelity increase of 50% with a maximum of 256%. Finally, when TimeStitch intelligently leverages periodic dynamical decoupling within its scheduling framework, a mean 64% improvement is observed over the baseline, relatively outperforming stand-alone dynamical decoupling by 19%, with a maximum of 287%.

# CCS Concepts: • Hardware → Quantum computation;

Additional Key Words and Phrases: Quantum computing, quantum information, quantum circuit optimization

## **ACM Reference format:**

Kaitlin N. Smith, Gokul Subramanian Ravi, Prakash Murali, Jonathan M. Baker, Nathan Earnest, Ali Javadi-Abhari, and Frederic T. Chong. 2022. TimeStitch: Exploiting Slack to Mitigate Decoherence in Quantum Circuits. *ACM Trans. Quantum Comput.* 4, 1, Article 8 (October 2022), 27 pages. https://doi.org/10.1145/3548778

#### 1 INTRODUCTION

Quantum computation is a revolutionary information processing model that leverages quantum mechanical phenomena to solve intractable problems. **Quantum computers (QCs)** evaluate quantum circuits, or programs, in a manner similar to classical computers, but quantum information's ability to leverage superposition, interference, and entanglement is projected to provide QCs significant advantage in cryptography [45], chemistry [24], optimization [30], and machine learning [6] applications.

Current QCs are prototype devices; they are less than 1,000 qubits in size, and they do not implement fault-tolerant, error-correcting codes. These devices suffer from high error rates; noise is introduced during state initialization, gate application, and measurement procedures. In addition to errors during active operations, qubits are also vulnerable to noise during periods of inactivity. In today's quantum systems, especially superconducting devices, *decoherence* error in idling qubits causes state to degrade exponentially over time from phase accumulation and amplitude damping. Several near-term applications require a critical path, or longest path of dependent circuit operations, that is proportional to its circuit size [12]. Large qubit idle windows result from circuit compilation, introducing periods during computation that are highly susceptible to decoherence errors.

Decoherence-related noise in today's quantum systems is incredibly difficult to fully model. Much of this noise varies among devices as well as over time [49]. As a result, error mitigation with techniques that assume a more complete model of system error, such as Hahn spin-echo and dynamical decoupling, are challenging to apply in a standardized, one-size-fits-all approach on today's QCs. More optimum error mitigation will result with customized methods that take the unique needs of a quantum system and workload into consideration. Once the characterization of quantum error is improved, error mitigation can be enhanced during algorithm execution. Our work's fundamental goal is to mitigate both phase accumulation and amplitude damping unique to an algorithm and machine pairing without extending the duration of the original compiled circuit. We present a novel technique to optimize circuits by taking advantage of flexible scheduling within *slack windows*, or periods of qubit idling before its next operation. The benefits of our approach are achieved without extending circuit runtime through either increasing the circuit depth or introducing circuit partitioning.

Qubit slack may appear trivial in unmapped circuits, but the impact and duration of idling qubits becomes obvious post compilation. Many near-term devices, such as superconducting circuits, feature nearest-neighbor topologies with sparse connectivity across qubits. Unfortunately, many QC

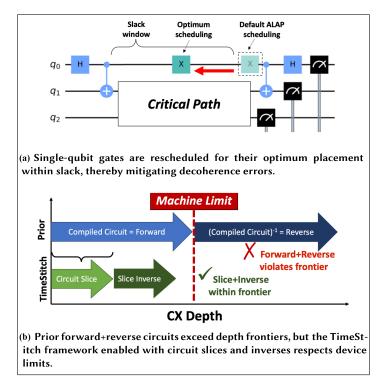


Fig. 1. Overview of the TimeStitch proposal.

applications have communication requirements that do not align well with these hardware constraints, resulting in the insertion of SWAP networks for intra-chip qubit communication. These SWAP networks increase the duration of quantum circuits, forcing a significant portion of physical qubit runtime, or time from state initialization until final measurement, to be suspended within slack.

Typical circuit scheduling methods such as "As Late As Possible" (ALAP) scheduling, a default approach in IBM Qiskit [4], assume that single-qubit operations are best placed at the end of slack windows. An example of ALAP scheduling is pictured in Figure 1(a) as an X gate in a dashed box. Our proposal reschedules single-qubit gates potentially away from its ALAP default position to an optimum placement within the slack window, also shown in Figure 1(a). The chosen gate placement is deemed optimum by measuring the fidelity at various gate positions in its slack window with a carefully designed tuning circuit. By choosing the optimal gate position, TimeStitch minimizes the impact of decoherence on qubits caused by dephasing and amplitude damping.

TimeStitch tuning procedures are implemented by intelligently leveraging the reversible nature of quantum computation. First, a tuning circuit starts with a slice of the original circuit up to a target slack window. Next, a window equal to the slack found in the original circuit is placed in the tuning circuit. Finally, the circuit slice is inverted to "undo" previous computation, returning all qubits to their original input state. The approach is thus referred to as a "slice + inverse" (SI) technique. Critical to the near-term, TimeStitch tuning employs reversibility without exceeding the machine fidelity limits on circuit depth. Prior work exploiting reversibility for predicting circuit outcomes builds a concatenated "forward+reverse" of a quantum circuit in its entirety, resulting in double the depth of the original circuit [35]. Under the reasonable assumption that target circuits

8:4 K. Smith et al.

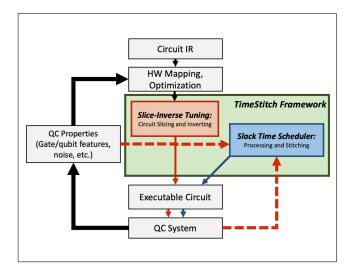


Fig. 2. QC compilation integrated with TimeStitch.

are already at or just under machine capacity in terms of critical path length, it is possible that the depth of such forward+reverse circuits can far exceed QC capability, which we refer to as the QC frontier. Here we assume that devices are paired with applications such that they are maximally utilized, and the QC frontier measures the number of *CX* gates in the critical path. Overstepping the QC frontier with optimization that employs quantum reversibility is shown in Figure 1(b) in blue; Section 2.3 discusses prior work that uses reversibility in greater detail. Alternatively, our proposal is constrained to specifically target slack windows whose corresponding SI, "slice + inverse," circuits are within the machine limits of in terms of circuit depth. This is shown in green in Figure 1(b). Further, we show that even with this constraint we are able to reap most of the potential benefit from our decoherence mitigation approach, because it still allows us to target most of the larger slack windows that both have higher potential for gate position tuning based benefits and create circuit slices of lower gate depth due to the very definition of slack.

Figure 2 depicts a quantum compilation flow that includes the **TimeStitch (TS)** slack optimizing scheduler that consists of two components. The first component is a module that identifies slack windows within a compiled quantum program and develops slack tuning circuits through circuit slicing and inversion procedures. These circuits are then used to independently optimize single-qubit gate positions within individual slack windows. Compilations and device properties are subject to variation, thus the optimum placements determined by slack tuning will be unique for each circuit and machine at a given period of time. TimeStitch is an empirical approach that does not attempt to fully understand the noise of a system with all its complexities. Instead, it forms a snapshot of the impact of noise during slack periods in quantum circuits during Slice-Inverse tuning. This information is used to determine optimum scheduling that mitigates error in each window, and TimeStitch creates a final "stitched" executable that is an optimized form of the original compiled circuit.

We note that our proposal to exploit slack in quantum circuits is not limited to our primary error mitigation approach of tuning single-qubit gate positions within slack windows. As a secondary contribution, we show that other techniques for error mitigation, such as dynamical decoupling [56], can also be specifically targeted within these slack windows through our generalized compilation framework. We show that when gate scheduling is intelligently coupled with periodic dynamical decoupling within the TimeStitch framework, the error mitigation techniques

compliment each other, resulting in even greater fidelity improvements across a variety of quantum circuits. Finally, the hallmark theme of exploiting slack in quantum circuits has significant parallels to slack-based optimization in classical computing, such as those at the circuit level as well as the microarchitecture level; we dive deeper into these parallels in Section 7.

To summarize, this work makes the following contributions:

- We observe the creation of slack windows as a result of compilation. To the best of our knowledge, we are the first to identify their potential for optimal quantum gate scheduling.
   We develop a framework that optimally schedules single-qubit gates for mitigating decoherence error that degrades idling qubit state.
- To the best of our knowledge, we are the first to exploit quantum reversibility toward gate scheduling, and importantly, in a manner cognizant to device depth limitations. Reversibility enables the mitigation technique to adapt to the unique characteristics across both applications and QCs to provide a solution that is not "one-size-fits-all."
- We design a slack analysis and circuit construction method that analyzes compiled QCs, identifies slack windows, and "slices" the original circuit to isolate dependency graphs up until instances of slack. These "slices" are then combined with a delay line equal to the corresponding slack window followed by the slice inverse circuit to create a total circuit that evaluates to a ground truth: the slice input state.
- We design and implement the TimeStitch Slice+Inverse (TS-SI) slack time scheduler that optimizes the scheduling of single-qubit gates within slack. Local optimals are learned during tuning procedures when individual slack windows are searched within slice+inverse circuits (above) to maximize the fidelity of the trivially known ground truth. TS-SI then "stitches" a final quantum circuit with optimum placements identified from tuning. During tuning procedures and final circuit creation, the bounds of the original circuit depth are respected as criteria for tuning (TimeStitch with Circuit Slice+Inverse Tuning plus Criteria (TS-SI+C)).
- We implement TimeStitch to suit deployment on real quantum machines and offer insights that can improve the realistic design of future quantum optimization proposals. The framework is evaluated on a variety of benchmark circuits transformed by baseline compilation and TimeStitch Slice+Inverse rescheduling. We compare TimeStitch against other scheduling heuristics such as ALAP, "as soon as possible" (ASAP), and Middle, all discussed in Section 5.4, and highlight TimeStitch's greater benefits over "one-size-fits-all" gate scheduling solutions.
- We show that our general TimeStitch compilation framework for targeting slack windows
  can encompass additional error mitigation techniques like periodic dynamic decoupling
  (DD) during schedule optimization. Analysis is provided to show that the two approaches
  can harmonize to create highly optimized circuits. Our highest performance gains derive
  from TimeStitch invoked synergistically with DD that implements our empirically derived
  DD heuristic. This implementation of TimeStitch significantly outperforms a standard DD
  approach.

TimeStitch holds great potential for impact in the area of quantum compiler design as it is the first proposal to exploit optimum scheduling of quantum operators within slack windows. While many existing techniques for mitigating error rely on adding extra gates to the circuit [13, 20, 56] calibrating new gates [50], or extending a circuit's runtime [32], TimeStitch leverages the gates already present in a quantum program in its base form. TimeStitch, however, can be invoked with DD optimization to reap the combined benefits of multiple state-of-the-art decoherence mitigation techniques. Additionally, a novel aspect of our framework is that unlike previous proposals

8:6 K. Smith et al.

that employ reversibility through "forward" and "reverse" circuits [35], program duration is not extended either during tuning procedures or in the final rescheduled circuit. This is critical in the near-term where QCs are aggressively pushed to the brink in terms of utilization.

This article proceeds as follows: Section 2 presents background information describing fundamental elements of this study. Section 3 details theory related to quantum computing and quantum error mitigation that motivates TimeStitch. Section 4 describes the design of the TimeStitch framework. Section 5 includes the methodology surrounding TimeStitch development and evaluation. Section 6 evaluates TimeStitch with experiments performed on real QCs. Section 7 is a discussion of future directions for TimeStitch as well as related work in both the areas of quantum circuit optimization and classical slack exploitation. Section 8 offers conclusions.

#### 2 BACKGROUND

## 2.1 Quantum Information and Near-Term QCs

The basic unit of quantum information is the quantum bit, or qubit. Qubits, unlike classical bits that hold a static value of either 0 or 1, demonstrate superposition in the form of  $|\psi\rangle = \alpha |0\rangle + \beta |1\rangle$  where probability amplitudes  $\alpha, \beta \in \mathbb{C}$  hold values such that  $|\alpha|^2 + |\beta|^2 = 1$ . Upon measurement,  $|\psi\rangle$  collapses into *either*  $|0\rangle$  or  $|1\rangle$ , effectively becoming a classical bit. A system of n qubits requires  $2^n$  amplitudes to describe the state.

Before measurement, qubits are manipulated with operations, or gates, to modify the quantum state's probability amplitudes. Quantum operations are unitary, and as a result, they are characterized as reversible with the same number of inputs as outputs. Unlike classical computation, there are many non-trivial single-qubit gates such as  $R_x(\theta)$  and  $R_z(\phi)$  that rotate the state around the x-and z-axes, respectively. An example of  $R_x(\theta=\pi/2)$  and  $R_z(\phi=\pi/2)$  rotation of qubit visualized with the Bloch Sphere is pictured in Figure 3(a). Pairs of qubits can be manipulated via multi-qubit interactions. One of the most common of these gates is the two-qubit, controlled- $(R_x(\pi)=X)$ , or CX gate. Together with single qubit gates, CX enables universal quantum computation. There are many choices of basis gate sets specified by the underlying hardware. For more information on the fundamentals of quantum computation we refer to Reference [34].

Current QCs, sometimes called **Noisy Intermediate Scale Quantum (NISQ)** devices, are error prone and less than 1,000 qubits in size [38]. These devices are extremely fragile, and as a result, some of the biggest challenges that limit scaling include limited coherence, gate errors, readout errors, and connectivity. Error during computation can corrupt results, but once noise is identified and characterized, it can often be effectively mitigated or corrected with software. Primary causes of loss of performance include decoherence and crosstalk.

Many errors in quantum systems arise from environmental coupling. For example, amplitude damping describes the sporadic loss of energy resulting in the  $|1\rangle$  state falling to the  $|0\rangle$  state; the rate of this process is described by the device's  $T_1$  time. Similarly dephasing, also referred to as phase accumulation or phase damping, details the sporadic change in relative phase and is expressed by the  $T_2$  time of the qubit. Both cause qubit state decoherence. Finally, crosstalk refers to error caused by simultaneous execution of gates on nearby qubits. The severity of each type of noise varies per qubit and calibration cycle.

We propose TimeStitch for the mitigation of decoherence errors. This is achieved by tuning single-qubit gate positions within idle periods in circuits, Section 3.

# 2.2 Qubit Idling in Compiled Circuits

Qubit runtime during circuit execution is the period spanning the first gate after qubit reset up until measurement. During its runtime, a qubit will spend some cycles in computation and others

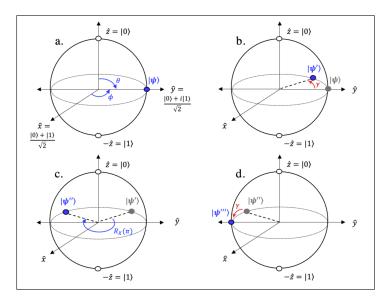


Fig. 3. Phase accumulation mitigation through Hahn spin-echo techniques. (a) A qubit  $|\psi\rangle$ , prepared with  $R_X(\theta=\pi/2)$  and  $R_Z(\phi=\pi/2)$ , rests on the y-axis of the Bloch sphere. (b) As time elapses, the phase of  $|\psi\rangle$  decays, and noise in the form of  $R_Z(\gamma)$  creates the quantum state  $|\psi'\rangle$  after a counterclockwise rotation around the z-axis. (c) A  $R_X(\pi)$  is applied to the qubit to produce  $|\psi''\rangle$ , and (d) the effects of dephasing begin to constructively interfere with  $|\psi''\rangle$  to produce the phase-coherent state  $|\psi'''\rangle$ . Another  $R_X(\pi)$  pulse restores  $|\psi\rangle$  from  $|\psi'''\rangle$ .

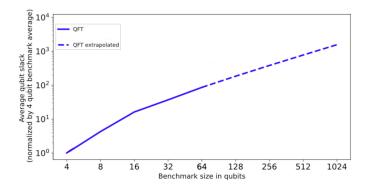


Fig. 4. Average qubit slack normalized by smallest, 4-qubit implementation vs. benchmark size in qubits for QFT benchmark. QFT circuits with 4-64 qubits are mapped to 65-qubit IBM QC, and slack trends are extrapolated to 1,024 qubits in anticipation of near-term machines.

idle as neighboring qubits undergo state evolution caused by gates on a circuit critical path. Idle time is referred to as slack.

Limited connectivity in near-term devices requires *SWAP* networks for qubit communication in mapped circuits. As we move toward larger devices, connectivity is anticipated to stay low as architectures such as heavy-hex topologies are expected to be the most favorable to scale superconducting qubit machines [9, 33]. As demonstrated in Figure 4, slack within circuits increases with the number of qubits because of limited qubit-qubit communication. In this plot, the **Quantum Fourier Transform (QFT)** is mapped to the 65-qubit IBM Q Manhattan quantum machine for

8:8 K. Smith et al.

4-, 8-, 16-, 32-, and 64-qubit implementations. Maximum optimization is used by the IBM Qiskit [4] transpiler. The slack windows that appear in each compiled circuit after the qubit runtime begins are identified, and the total time idling within circuit slack is averaged between all qubits. This average qubit slack for all implementations is then normalized by the smallest slack average corresponding to the four-qubit QFT instance. Figure 4 includes a plot (solid line) of the normalized, average qubit idle time total for the 4- to 64-qubit QFT circuits mapped to the 65 qubit QC. As it is anticipated that nearest-neighbor QCs will scale to thousands of qubits in the near-term, the line detailing average slack is extrapolated (dashed line) from 64 to 1,024 qubits to anticipate future technologies. The QFT extrapolated trends show that the circuits have average slack that increases by factors of approximately 1,000× at 1,024 qubits, demonstrating that the amount of qubit inactivity during its runtime has a direct relationship with circuit size when mapped to near-term hardware. Many quantum algorithms are anticipated to demonstrate this same trend, and regardless of QC technology, QC applications will experience increased circuit slack as algorithms and critical paths scale without substantial parallelization.

By default, compilation tools tend to schedule single-qubit operations within slack windows for ALAP, meaning that gates will not execute until another operation, typically either a measurement or a two-qubit operation along a critical path, can occur immediately afterwards [4]. Scheduling qubit operations for ALAP assists with mitigating noise associated with  $T_1$  and  $T_2$  decoherence if qubit runtime has not initialized. ALAP execution, however, is not always ideal once a qubit holds state and is more vulnerable to decoherence. Rather than tolerate slack as an unavoidable artifact of compilation and assume ALAP gate defaults, we are motivated to explore theoretical and practical techniques for decoherence mitigation during the periods where qubits idle, as illustrated in Figure 1.

# 2.3 Considerations with Applied Reversibility

Quantum computation is reversible, because quantum operations are unitary. A requirement for a unitary operation, U, is that  $UU^{-1}=U^{-1}U=(ID)$  where  $U^{-1}$  is the operation inverse and (ID) is the identity operation. The identity operation does not evolve qubit state and produces an output equal to the input; it acts as a fixed-duration, "do nothing" instruction. As a note, quantum circuit measurement is not reversible as it collapses superimposed qubits into a classical bitstring.

A quantum circuit followed by its logical inverse, or a "forward+reverse" circuit, produces circuit, thus ideally produces the original or initial state. In QRAFT [35], quantum reversibility reduced error in circuits by increasing the likelihood of determining the correct evaluation output. Since the outputs are known as a ground truth for forward+reverse characterization circuits as they are equal to the initial state, noisy QC results can train a machine learning model to discern error attributes for a machine. The model is used to predict true quantum circuit outcomes when circuits are in their forward+reverse form.

While Reference [35] provides a boost in quantum circuit accuracy, it assumes the ability to successfully run quantum algorithms where critical paths, or depths, are twice that of the original circuit. This may be a reasonable approach for small quantum circuits that terminate well within the bounds of coherence times, but hardware is ideally maximally utilized in practical workloads. Thus, circuits operating at the boundary of machine thresholds may produce unreliable results if executed in their extended forward+reverse form. To avoid observing a noisy distribution, techniques invoking reversibility should consider the duration of the original quantum circuit as bounding criteria.

In this work, quantum reversibility is leveraged by TimeStitch to enable the optimization of single-qubit placement within slack. Unlike Reference [35], which applies reversibility toward predictive models, we utilize the true output provided by inverting quantum circuits to produce circuit

schedules that outperform baseline ALAP compilations. These improvements are achieved without exceeding the critical path criteria either during slack tuning or in the final, optimized circuit. A full description of the TimeStitch framework is found in Section 4 with details about circuit depth constraints in Section 4.3.

# 2.4 Spin-echo Error Mitigation: Dynamical Decoupling

To preserve quantum state without corrective codes, open-loop error mitigation can be applied to refocus signals. An example of this type of correction is DD [56] that "decouples" compute qubits from environmental noise. The most elementary form of DD suppresses single-qubit phase accumulation with Hahn spin-echo techniques where  $R_x(\pi) = X$  instructions are insert into circuits during runtime. These instructions reverse the impact that dephasing has on quantum state over time. For example, consider a quantum state  $|\psi\rangle = \frac{|0\rangle + i|1\rangle}{\sqrt{2}}$ . This qubit on the positive y-axis of the Bloch sphere is pictured in Figure 3(a). Ideally,  $|\psi\rangle$  would hold state information for infinite time, but phase information is highly susceptible to decoherence. In Figure 3(b), the decay of state by the unknown rotation  $R_z(\gamma)$  causes  $|\psi\rangle$  to evolve to  $|\psi'\rangle$ . Hahn spin-echo techniques apply a  $R_x(\pi)$  operation to  $|\psi'\rangle$  in Figure 3(c) to mitigate the phase accumulation caused by decoherence, resulting in state  $|\psi'\rangle$ . The continued dephasing shown in Figure 3(d) counteracts the original rotation of  $R_z(\gamma)$ , refocusing phase information to produce qubit  $|\psi'''\rangle$ . Restoring the original state  $|\psi\rangle$ , pictured in Figure 3(a), with phase information intact, is possible with the application of a final  $R_x(\pi)$  pulse to  $|\psi'''\rangle$ . The procedure of inserting  $R_x(\pi)R_x(\pi) = XX$  mid-circuit preserves the semantics of the original circuit as  $UU^{-1} = (ID)$  where (ID) is the identity operation.

Many different forms of DD have been proposed [39, 47, 54], and DD has shown promise on near-term quantum processors [5, 10, 13, 23, 37]. While DD has considerable potential, the quantum community is still far from widespread implementation due to limitations stemming from non-ideal properties and overheads of the decoupling pulses [26, 48]. In fact, past work demonstrated that naively implementing DD in a universal manner on idle qubits can result in decreased circuit fidelity [13]. Thus, there is still significant room for DD improvements, both stand-alone as well as combined with other error mitigation and correction techniques.

As DD is a leading technique for decoherence mitigation, determining its optimal use in conjunction with TimeStitch was worthy of exploration and resulted in considerable benefits. Section 3.3 includes more discussion motivating integration of DD into the TimeStitch framework.

### 3 THEORY FOR SLACK WINDOW OPTIMIZATION

## 3.1 Tuning Gate Positions for Phase and Amplitude Errors

DD techniques employ additional gates to recohere quantum state in the presence of noise. Rather than add gates to a circuit, we are motivated to search for ways to refocus signals using operations already present within the circuit. In the most simple example of how gate placement within slack could influence circuit outcomes, consider the case where a qubit in excited state,  $|\psi_{initial}\rangle = |1\rangle$ , enters an idle period. If the next gate acting on the qubit is an X gate that would NOT the state of the qubit, lowering it to the ground state,  $|\psi_{final}\rangle = |0\rangle$ , then the preferred execution schedule would be ASAP to avoid amplitude damping from negatively impacting  $|\psi_{initial}\rangle$  as it idles. Conversely, if  $|\psi_{initial}\rangle = |0\rangle$  at the beginning of slack, then there are advantages to scheduling an upcoming X gate for ALAP to extend the time the qubit spends in the ground state that is less susceptible to noise.

Quantum states are often more complex superpositions than those described in the aforementioned example. For this reason, the circuit in the top of Figure 5 is used as a micro-benchmark to demonstrate the viability for decoherence mitigation via gate rescheduling within slack. An

8:10 K. Smith et al.

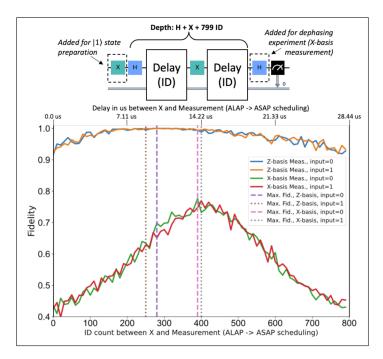


Fig. 5. Demonstration of amplitude damping and dephasing correction via Hahn spin-echo techniques. The pictured H+X+Delay circuit on a single qubit tunes X gate placement within a slack window to relate position to state fidelity. Measurement in the Z-basis ( $|0\rangle/|1\rangle$ ) captures amplitude information. An H at the circuit end causes an X-basis measurement ( $|+\rangle/|-\rangle$ ), capturing phase information. When X is scheduled near the middle of the slack window, the fidelity is maximized. Maximum location for each experiment differs.

IBM QC was used for this Hahn spin-echo inspired micro-benchmark experiment. The core of the circuit consists of an H gate that puts a qubit into superposition, a slack window artificially created with 799 identity (ID), or "do nothing," operations, and an X gate that is tuned within the slack. To tune X, the 799 (ID) gates are distributed between two partitions on either side of the X gate that can range from 0 to 799 (ID) gates in size as the X gate sweeps the slack. Additional components included in a select subset of micro-benchmarking experiments are an X gate that prepares the input state  $|1\rangle$  and an H before measurement that allows measurement to be in the X-basis ( $|+\rangle = \frac{|0\rangle + |1\rangle}{\sqrt{2}}$ ,  $|-\rangle = \frac{|0\rangle - |1\rangle}{\sqrt{2}}$ ) rather than in the Z-basis ( $|0\rangle$ ,  $|1\rangle$ ). These additions are shown in dashed boxes.

The circuit in Figure 5 is inspired by  $T_1$  and  $T_2$  experiments, but here we do not seek to measure decoherence times. Instead, each of the four versions of the micro-benchmark (input  $|0\rangle$ /measurement Z, input  $|1\rangle$ /measurement Z, input  $|0\rangle$ /measurement X, input  $|1\rangle$ /measurement X) has a fixed duration while the X gate position is tuned in search of a maximum fidelity schedule. As a note, each (ID) gate has a duration equal to that of a single X gate on the IBM QCs: approximately 35.56 ns. We define fidelity as the Hellinger fidelity between an ideal distribution and the distribution produced from a real QC run. The graph in the lower half of Figure 5 demonstrates that gate placement within slack can influence circuit outcome. Final measurement in the Z-basis ( $|0\rangle/|1\rangle$ ) captures information about amplitude damping. An H at the circuit end causes an X-basis measurement ( $|+\rangle/|-\rangle$ ), capturing information about qubit phase decoherence. When X is scheduled near the center of the slack window, the fidelity is maximized in all four circuits,

although the benefits associated with phase correction were more substantial. This result shows that even though we are not implementing true DD error mitigation, rescheduling *inspired* by Hahn spin-echo techniques can effectively correct both dephasing and amplitude damping error.

The maximum fidelity schedule for each experiment differs in Figure 5, suggesting the importance of state and measurement basis for optimum placement. In realistic workloads, many variables exist such as variation in single-qubit gate rotation, the qubit that the gate acts on, the slack window state, and the slack duration. The theory alone does not provide a clear prediction of optimum schedule for general use cases, motivating the need for automated solutions that rely on empirical observations, which we pursue by exploiting the quantum property of reversibility.

# 3.2 Understanding Real-Machine Impact

- 3.2.1 Crosstalk. Crosstalk is the accidental transfer of a qubit's information to surrounding qubits. Two adjacent gates, especially two-qubit interactions, executed simultaneously and within close proximity on nearest-neighbor QCs often experience lower gate fidelity as a result of crosstalk. Because of the severity of crosstalk, software mitigation techniques have been proposed [14, 32]. Studies have shown that single-qubit, single-qubit crosstalk is trivial [32]. Thus, the scheduling of single-qubit gates in adjacent slack windows can be tuned independent of one another. Discussion our framework's slack tuning procedures is found in Section 4.
- 3.2.2 Variation in Qubit Characteristics. Near-term quantum machines are affected by non-deterministic spatial and temporal variations in their characteristics. For instance, prior work [49] observed the prevalence of a wide distribution of machine characteristics with considerable spatial and temporal variation. From the spatial perspective, they observe the coefficient of variation, calculated as standard deviation over mean, to be in the range of 30-40% for  $T_1/T_2$  coherence times, as well as nearly 75% for two-qubit error rates. From a temporal perspective, they observe more than  $2\times$  variation in error rates in terms of day-to-day averages.
- 3.2.3 State Diversity within Slack Windows. Each quantum algorithm has a unique objective, resulting in a large amount of state variation during computation, especially within slack. As mentioned, every QC has a distinct noise signature with impact of varying severity depending on an idling qubit's state value. Unfortunately, certain states are more vulnerable to error. For instance,  $|1\rangle$  is more vulnerable to  $T_1$  amplitude dampening than  $|0\rangle$ , and  $T_2$  dephasing is highly influential to superimposed states such as  $\frac{|0\rangle+|1\rangle}{\sqrt{2}}$ .

Because of variation within quantum machines, see Section 3.2.2, and how this variation impacts circuits, it is challenging to develop an umbrella benchmark, or a set of benchmarks, for slack tuning that accurately captures unknown state and error attributes seen in real QC execution. Thus, we are motivated to use the circuits and the machines under investigation themselves, building upon the reversible nature of quantum computation, as the basis for slack tuning to accurately capture execution diversity while searching for optimum gate schedules.

# 3.3 Considerations for Invoking Dynamical Decoupling

The XX sequence implements an elementary form of DD that provides Hahn spin-echo correction of phase accumulation. DD, however, requires additional gates within the correction sequence, because rotation operations around at least two axes are necessary for more robust qubit error decoupling [47]. DD with a single "universal decoupling" sequence requires four gates:  $R_X(\pi)R_y(\pi)R_x(\pi)R_y(\pi) = XYXY$  [56]. The universal decoupling sequence adds increased protection to quantum state, because  $\pi$  rotations about both the x- and y-axes make the qubit more resilient to environmental noise. Additionally, Reference [53] analytically shows that XYXY is the

8:12 K. Smith et al.

superior choice for DD correction of arbitrary quantum states when considering DD sequences containing four gates on two axes.

DD has proven effective at correcting single-qubit states and, to a lesser extent, two-qubit entangled states in superconducting systems [37]. In addition, DD in the form of the XX sequence to implement Hahn-echo correction has also improved the Quantum Volume of a real QC in concurrent work [23]. Both of these demonstrations, however, cost additional circuit instructions during runtime. When inserting DD sequences into a circuit for signal refocusing, the number of additional gates should be carefully considered as gate errors tend to accumulate, potentially destroying the state of the system rather than protecting it from environmental impact [47]. Single-qubit gate errors on superconducting qubits are on average of order  $10^{-4}$  [1, 23], and although individually small, collective errors can degrade circuit performance, especially as circuits scale on maximally utilized machines.

The problem of diminishing quantum circuit outcomes with a naive, universal DD implementation is discussed in related work in a framework called ADAPT [13]. ADAPT proposes a clever idea, to evaluate potential DD insertion by transforming target circuits into "decoy" circuits containing just Clifford gates. These novel decoy circuits can then be tractably simulated and selective DD insertion strategies evaluated. This study shows that in general, there is not a one-size-fits-all solution for DD, but typically adding some DD to a circuit provides improvements. Although impressive performance gains for a small subset of benchmarks are reported using an evaluation metric based on "total variation distance," the benefits of Reference [13] may not be as substantial using more standard metrics such as Hellinger fidelity or probability of success. Additionally, the Clifford approximation used in ADAPT fails to fully model internal states of the circuit, so it is unlikely that the implemented DD is optimum unless the benchmark consists of mostly Clifford operators. In this article, we propose an alternative that uses circuit slicing and uncomputation to tune DD as well as gate location within paths with high slack using execution on the actual machine. Our "slice+inverse" approach avoids the inaccuracy of the Clifford approximations but will have some tradeoffs of execution overhead and critical path limitations, as discussed in Sections 5.3 and 4.3.

In addition to selecting the proper gate sequences and locations for DD within circuits, timing must also be considered so that DD effectively "unwinds" error on decaying quantum states. In other words, there must be enough spacing between the execution of gates in the DD sequence to provide corrective benefits [7, 23]. A commonly implemented form of DD with uniformly spaced correction gates, such as with *XYXY*, is referred to as periodic DD [7].

For effective application of DD to circuit optimization, it must be deployed through intelligent tuning routines that avoid scenarios of introducing additional gate error that outweighs corrective benefits. As discussed later in Section 4.4, intelligent tuning can be especially beneficial when DD is deployed in conjunction with other error mitigation techniques. The TimeStitch compilation framework is expanded to incorporate additional decoherence mitigation in the form of periodic DD using XYXY, and we empirically tune DD parameters for maximum overall benefit. Empirical details of the methodology are located in Section 5.4.

#### 4 DESIGNING THE TIMESTITCH FRAMEWORK

## 4.1 Lessons from the Theory

Section 3 motivates the need for empirical solutions that are efficient in utilizing the quantum machines. To do this, these approaches should ideally be backed by robust quantum theory that also take into consideration the abilities of near-term machines. For slack optimization, our proposal for a practical approach is built on the following theoretical lessons.

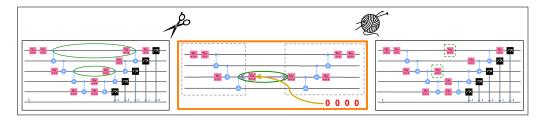


Fig. 6. TimeStitch Framework. Left: Circuit is compiled to the target machine and slack windows for tuning are identified throughout the quantum circuit. Middle: In the absence of known circuit outcomes, gate positions are optimized by exploiting quantum reversibility. For each slack window, a circuit slice from circuit start to the slack boundary is constructed and concatenated with a delay line and its inverse. The gate position in the target window is tuned with the goal to make circuit output match the input ( $|00...0\rangle$ ), implying position of maximum fidelity. Right: Tuned gate positions are stitched together to construct an optimized circuit schedule.

- ① *Prevalence of slack windows*: Section 2.2 describes that slack windows exist in compiled quantum circuits, and their amount and duration are correlated to the size of the quantum circuit targeted for a QC demonstrating limited connectivity between physical qubits.
- ② Adjusting gate positions: Opportunities for improving the fidelity of quantum circuits exist through adjusting the execution of single-qubit gates from ALAP scheduling to earlier placement within slack windows. A proof-of-concept case study using a micro-benchmark is presented in Section 3.1.
- ③ Optimal positioning: Optimal gate scheduling within a slack window depends on gate and qubit characteristics along with input qubit state. The vast space of these parameters on real machines, Section 3.2, suggests that offline machine characterization on test inputs and circuits is insufficient and impractical for finding optimal gate positions for general use cases.
- ④ Reversibility: Although we cannot predict the outcome of a quantum circuit execution, Section 2.3 describes that quantum reversibility can be used to provide a ground truth. We are motivated to apply reversibility to learn properties of quantum circuits and machines within windows of slack to implement application-specific decoherence mitigation.
- ⑤ *Impact of single-qubit crosstalk*: Minimal impact of single-qubit crosstalk, Section 3.2.1, means that the single-qubit gate position in each slack window can be optimally tuned independent of those in other slack windows.
- ⑥ Synergistic TS deployment with prior art: Decoherence mitigation techniques such as DD exist and are shown to be effective, but these solutions must be carefully implemented so that the operational characteristics unique to a circuit and machine pairing are considered. In the context of this work, a slack window must be of a minimum duration to provide adequate spacing between DD sequence gates within the window so that DD is effective and DD gate errors are trivial. Integrating the proposed TS technique with DD involves re-evaluating the best windows to incorporate DD as slack windows are divided when gate positions are adjusted. Details are found in Sections 2.4 and 3.3.

# 4.2 A Practical "Slice + Inverse" Approach

The practical TimeStitch approach leverages the quantum phenomenon of reversibility to adjust the execution timing for single-qubit gates within slack windows through the process of circuit SI. An overview of the framework is shown in Figure 6 and is discussed below.

8:14 K. Smith et al.

4.2.1 Baseline Compilation of the Quantum Circuit. The TimeStitch framework begins with a quantum circuit compiled from a device-independent intermediate representation (IR) into machine-ready code. The baseline circuit, methodology discussed in Section 5.3, appears in the left circuit of Figure 6.

- 4.2.2 Identifying Slack Windows. The TimeStitch framework identifies quantum circuit slack windows after baseline compilation. The identification procedure requires traversing the components of a quantum circuit that implements default ALAP scheduling from end to end. During this procedure, slack windows are found, and their durations are calculated using gate timing data collected from the QC. A subset of windows are identified that contain single-qubit operators eligible for rescheduling within slack. Two such windows are circled in green in the left circuit of Figure 6. As a note, we do not consider the time before the first operation on a qubit as slack, since the qubit is uninitialized and its runtime has not begun.
- 4.2.3 Generation of Slice+Inverse Calibration Circuits. From the identified slack, calibration circuits consisting of sliced partitions of the original circuit and their corresponding inverses are generated. Tuning experiments determine optimal schedules for single-qubit gates that are suitable candidates depending on criteria set by the depth of the original circuit targeted for optimization. Setting a criteria for tuning circuit depth is critical for ensuring that calibration procedures do not exceed the frontier of the QC.

We restate that our goal is to find the optimum gate position within each slack window to boost overall circuit fidelity. We employ the property of quantum reversibility, described in Section 2.3, to determine optimum single-qubit execution times within circuit slack to mitigate decoherence and thus maximize the fidelity at the end of the slack.

For each eligible slack window, for example, the window circled in green in the center circuit of Figure 6, a circuit slice is constructed, terminating at the end point of the particular window. Implementation-wise, this circuit slice is simply a subcircuit of the original circuit consisting of the dependency graph up until the end of the slack window. Qubit mapping of the original circuit is preserved in the subcircuit, thus the output of the circuit slice emulates the activity of the circuit up to the end-point of the slack window. The inverse of the circuit slice is then constructed and concatenated at the endpoint of the slice. An example slice and its inverse is shown in dashed boxes in the center circuit Figure 6; we refer to this as an SI circuit. Measurement operations, not shown in Figure 6, are added to the end of the SI circuit. In an ideal, noise-free setting, the concatenation of a slice and inverse would produce the slice's input as the output of the inverse because of quantum reversibility. In a realistic, noisy setting, our goal is then transformed to tuning the gate position in the slack window so that the probability of achieving the slice input state as the output of the total concatenated circuit is maximized. This is equivalent to maximizing the circuit fidelity of the original slice under the reasonable assumption that noise impact on the slice and its inverse are well correlated.

The input state to the slice is trivially known if the slice is constructed from the start of the entire quantum circuit; it is the ground state or  $|00\ldots0\rangle$ . As a result, the target output of the SI circuit is also the ground state  $|00\ldots0\rangle$  as shown in red in the center of Figure 6. Since input states, gates, and noise characteristics all influence the optimal gate position, each slack window must be sliced individually. TimeStitch creates the required, unique SI circuits in an automated manner, and the total number of SI tuning circuits for an input circuit is equal to the number of identified slack windows with single-qubit operations.

As a note, the maximum depth of an SI circuit is approximately twice the depth of the original circuit if a slice extends to near the circuit's end point. In the approach discussed in this Section, we do not limit the depth of the SI circuits; constraining the depth is discussed in Section 4.3.

4.2.4 Optimal Gate Placement in Slack Windows. Optimal gate placement within slack windows is determined by locating the position within each SI slack window where the ground truth  $|00...0\rangle$  is maximized. A variety of search strategies can be employed to find the local slack optimum, but optimizing the window search is orthogonal to this work. Each window is optimized independently with its corresponding SI circuit, and resolution of the search can be selected based on the availability of quantum machine resources. In other words, although tuning overhead is manageable and worthwhile for notable quantum circuit fidelity gains on real hardware, if a user is limited by the number of available quantum circuit runs, then each calibration circuit can be more coarsely searched.

4.2.5 Stitching the Per-Window Positions Together. After the optimal schedules are estimated via SI tuning, the local schedules are stitched together to form a composite, rescheduled circuit, as pictured in the right of Figure 6.

# 4.3 Constraining by Circuit Depth

The total number of SI tuning circuits is equal to the number of slack instances that contain tunable single-qubit gates. However, some of these SI circuits, such as those that cover slack appearing at the end of the input circuit, may have a depth exceeding that of the original circuit. This is because the SI tuning circuit will have a depth twice that of the subcircuit slice leading up to the slack window, as seen in the center of Figure 6. Past work [35] leveraging reversibility does not take circuit depth increase into consideration, but not doing so could potentially push beyond the frontier of the targeted device. The intuitive reasoning is that in the near term we are likely to be executing quantum applications that are already at the brink of a QC's capability in terms of the machine's critical circuit depth. Building tuning circuits beyond this critical depth can be detrimental to optimizing the original circuit, because it may provide false optimum schedules that are distorted by noise. To be mindful of the limitations of gate error and decoherence in current QC hardware, TS-SI can be initialized using the depth, or critical path, of the original circuit as bounding criteria. This version of TimeStitch, illustrated earlier in Figure 1(b), is known as TS-SI+C.

Here we focus only on the two-qubit operations along the critical path as a measure of depth. Two-qubit operations dominate in terms of influence on program output, because, on average, their error rates and duration are  $>10\times$  of a single-qubit operation [23]. This is particularly favorable for TimeStitch, because large slack windows have two-qubit depth that are often considerably lower than the critical depth, which is why large slack windows exists. Moreover, large slack windows are likely to provide substantial benefits due to the wider space for gate position tuning.

With TS-SI+C, depth is calculated for each of the SI tuning circuits. Those having a depth less than or equal to the depth of the original circuit are marked for use during TS-SI+C slack window gate position tuning. All untuned slack windows maintain default ALAP scheduling. Examples of circuit locations eligible and ineligible for TS-SI and TS-SI+C tuning with TimeStitch optimization are pictured in Figure 7. In the original compiled circuit used as TimeStitch input, slack windows 1, 2, and 3 are eligible for TS-SI as they all have tunable single-qubit operations. However, only slack windows 1 and 2 satisfy the depth criteria, since their corresponding SI circuits are of lower depth than the original circuit. Thus, only slack windows 1 and 2 are tuned by TS-SI+C. There are many other locations in the circuit, such as slack windows without single-qubit gates or periods before qubit runtime begins that are ineligible for slack tuning. This is also illustrated in Figure 7.

# 4.4 Integrating TS with Dynamical Decoupling

As mentioned in Section 3.3, DD sequence gates are spread within an idle window with adequate spacing between gates to provide maximal decoherence mitigation. Additionally, we wish to

8:16 K. Smith et al.

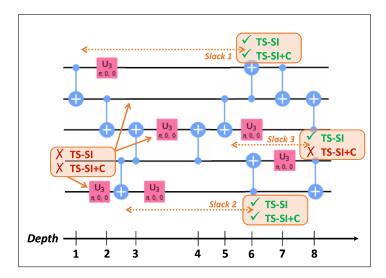


Fig. 7. Eligible and ineligible circuit locations for TS-SI and TS-SI+C tuning. Slack windows 1, 2, and 3 are eligible for TS-SI as they all have tunable single-qubit operations. However, only windows 1 and 2 satisfy the depth criteria and are thus tuned by TS-SI+C.

provide maximal correction benefits with DD without increasing the number of gates to the point where gate errors accumulate and degrade the state of the system [47]. Slack duration in compiled quantum circuits are prone to a vast amount of variation across applications and QCs, and implementing TimeStitch optimization results in the challenge of cutting large idle windows into smaller segments of size that are unknown before SI tuning procedures. DD implementation is highly dependent on window size, so this presents a challenge for employing DD to generate a maximally TimeStitch-scheduled circuit. Thus, we are motivated by the potential of DD to empirically develop a heuristic that generalizes DD to a vast set of use cases. This heuristic will be based on the periodic *XYXY* sequence and serves as an additional optimization benefit of the TimeStitch framework. More information is found in Section 5.4.

#### 5 METHODOLOGY

# 5.1 Evaluation Effort: Quantities and Constraints

We perform all our experiments on actual IBM superconducting quantum machines to faithfully capture true device characteristics. Our evaluations encompass roughly 2,500 quantum jobs to the cloud, comprising of over 600,000 circuits, with confidence built on a total of over 4,000,000,000 QC executions. Of these, we show results for 10 applications, each paired to a target machine satisfying the following criteria: (a) consistent machine availability, (b) non-negligible probability of the correct application output during baseline evaluation, (c) limited variability in correct output probabilities, and (d) maximized machine qubit utilization while respecting the previous constraints.

#### 5.2 Circuits for Evaluation

Our benchmarks are representative of real-world use cases, described here and in Table 1. Due to limitaions on circuit width because of machine size and depth because of coherence times on available near-term QCs, benchmarks that included six qubits or fewer and of shorter duration

Bench	Q	D	Output	# SW/ Cons.	Avg. SW (1e-6 s)	Avg. SW (1q count)	Dev
QAOA	4	15	0101 + 1010	10/5	0.85	23.90	Guad
QAOA	6	84	101000 + 111101	19/10	0.91	25.59	Jak
Gibbs	5	12	0000 + 0101 + 1010 + 1111	3/2	1.17	32.90	Jak
QFT	4	29	1010	15/8	1.25	35.15	Tor
GHZ	5	8	10101	3/3	1.48	41.62	Syd
VQE	4	63	0111	18/10	1.67	46.96	Guad
QFT	5	39	00101	25/12	1.83	51.46	Tor
QEC	5	26	00000 + 01011	4/4	2.41	67.77	Cas
Adder	6	64	000110 (1+0)	64/8	3.02	84.93	Syd
VQE	6	51	111111	13/5	3.99	112.20	Cas

Table 1. Benchmarks and Their Characteristics Sorted by Avg. SW

Q: number of application qubits; D: circuit depth in CXs; Output: application outputs; # SW/Cons.: number of slack windows/slack windows targeted under depth constraint; Avg. SW: average window wize in weconds and in count of single-qubit gates of duration 35.56 ns; Dev: target machine.

were included in TimeStitch evaluation. Brief descriptions of the benchmarks used in our study are as follows.

- *QFT*: QFT is a circuit used as a building block for applications such as Shor's algorithm for quantum factoring [45] and phase estimation. It converts a quantum state from the computational basis to the Fourier basis through the introduction of phase. QFT was constructed for five and five qubits [34].
- Quantum Approximate Optimization Algorithm: Quantum Approximate Optimization
  Algorithm (QAOA) [17] is a variational quantum-classical algorithm to solve combinatorial optimization problems. QAOA is implemented atop a parameterized circuit called an
  ansatz. We use one instance of a hardware efficient QAOA ansatz, and its solution is simple
  to predict when solving MAXCUT on a "ring of disagrees" graph structure. We use QAOA
  ansatz constructed for four and six qubits.
- Variational Quantum Eigensolver: Variational Quantum Eigensolver (VQE) [36] is a hybrid algorithm like QAOA and is used to variationally find the lowest eigenvalue of a given problem matrix by computing a difficult cost function on the QPU and feeding this value into an optimization routine running on a CPU. We implement VQE on a hardware-efficient SU2 ansatz [2] and use one instance as the benchmark. We construct the ansatz for four qubits and six qubits.
- Gibbs State Preparation: The preparation of Gibbs state has applications in quantum simulation, optimization, and machine learning. We take a VarQITE ansatz based approach to create the Gibbs state [59]. We use five qubits for the Gibbs circuit.
- Quantum Repetition Code Encoder: Error correcting codes are the means by which fault-tolerant quantum computers are able to execute arbitrarily long programs. Many such codes have been developed that make multiple tradeoffs [8, 15, 19]. Here, we target a error correcting repetition code encoder whose effect is to distribute the quantum information in the initial state across an entangled N-party logical state. This introduces redundancy to the encoding that can be exploited for error detection [44]. We use an encoder targeting five qubits.

8:18 K. Smith et al.

• Greenberger-Horne-Zeilinger State Preparation: Greenberger-Horne-Zeilinger (GHZ) state [21] generation is a non-traditional benchmark but useful as many complex quantum algorithms begin by entangling all qubits before computation in a state preparation process. In this benchmark, all qubits are first fully entangled before X gates swap the  $|0\rangle$  and  $|1\rangle$  probability amplitudes. Finally, qubits are unentangled to restore the input state. GHZ was implemented for five qubits.

• Ripple Carry Adder: Adders are a critical logic building block for quantum logic such as in Shor's algorithm for quantum factoring [45]. We implemented a linear-depth, two-bit ripple-carry adder quantum circuit that uses six qubits [12].

#### 5.3 Infrastructure

TimeStitch is implemented as a compilation pass that performs schedule optimization on top of a highest-baseline compilation of Qiskit Terra 0.16.4 to map and optimize for the IBM machines [4]. We distribute across five quantum devices: Casablanca (7 qubits), Jakarta (7 qubits), Guadalupe (16 qubits), Toronto (27 qubits), and Sydney (27 qubits). Machine details are on the IBM Quantum Systems page [3].

The IBM QCs are accessed via the quantum cloud. Resources are shared among hundrerds of thousands of users running more than two billion experiments per day [27], causing the queue time to service a quantum experiment request to vary significantly. As a result, an efficient means to utilize the cloud QCs is to maximize batches, or jobs, sent to a IBM QC. Jobs are treated as a single task, allowing for for the sequential processing and combined results of multiple circuits. A single quantum job of our target QCs can execute a batch of up to 900 circuits.

To keep the tuning overhead manageable, restricting calibration within a single job is key as the job runtime is more predictable and often significantly shorter than queue time in the IBM quantum cloud [43]. We use utilize this entire batch across the tuning of gate positions for different slack windows. Thus each slack window gets  $N = \frac{900}{\#SW}$  circuit slots for tuning, and the resolution of each window's gate position sweep is  $R = \frac{N}{SW_{length}}$ .

The benefits of our proposal on these circuits is evaluated on the **Probability of Success (POS)** metric that is the ratio of a number of error-free trials to the total number of trials—a common metric for evaluating quantum optimization.

## 5.4 Evaluation Comparisons

To analyze the effectiveness of TimeStitch, its performance was compared to other universal gate scheduling techniques as well as alternative error mitigation techniques targeted toward slack. The set used in comparision to TimeStitch are described below.

- *ALAP*: ALAP is the default scheduling technique implemented in the Qiskit compiler. All gates appearing in slack windows will be executed at the end immediately before the next two-qubit gate that acts as the slack end boundary. As ALAP is the default compilation setting of the Qiskit compiler, it acts as the baseline comparing the TimeStitch framework to other scheduling and optimization techniques described in this Section.
- *ASAP*: ASAP forces all gates appearing in slack windows to be executed immediately after the two-qubit gate that acts as the slack beginning boundary.
- *Middle:* Middle scheduling is a naive scheduling technique that executes all single qubits within slack at the center of their slack windows.
- *TS-SI:* TS-SI corresponds to the design from Section 4.2. The TS-SI form of the proposed optimization framework is "unconstrained" and invokes slack tuning on all circuit windows

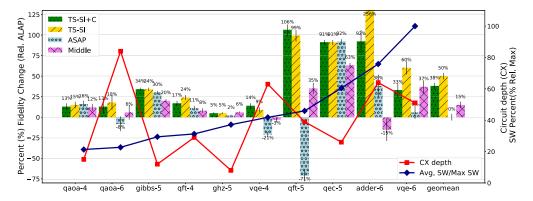


Fig. 8. POS benefits of different approaches over the ALAP baseline. TS-SI benefits are highest at 50% mean improvement over the baseline. TS-SI+C provides a 38% improvement. ASAP and Middle provide negligible and lower benefits on average, respectively, with detrimental individual outcomes on some individual benchmarks. In general, TS benefits increase as slack windows grow. These results were generated with real QC experiments.

using slice+inverse calibration circuits, regardless of circuit depth. The final optimized output circuit, however, is of depth equal to that of the original circuit.

- *TS-SI+C*: TS-SI+C corresponds to the design from Section 4.3. TS-SI+C is a practical approach as it respects the critical path of the original circuit. Windows are not tuned if their slice+inverse circuit exceeds the depth of the original circuit. Windows that are not tuned because of depth criteria violation maintain default ALAP scheduling.
- *DD*: In this article's implementation of DD, universal *XYXY* decoupling appears in slack windows. A single round is added to the window if it fits within the window duration. The gates are evenly spread across the window to create a periodic sequence.
- Dynamic Decoupling w/Heuristic: Dynamic Decoupling w/Heuristic (DD(H)) is similar to above, but DD is only added if a heuristic threshold inspired by References [23, 47, 53] is met to allow for adequate spacing between DD pulses. The inspiration for our DD heuristic is discussed further in Section 3.3. For our evaluation, the empirically found duration threshold is that the slack duration should be greater than or equal to four times the DD sequence duration.
- Integrated TimeStitch and Dynamic Decoupling: Integrated TimeStitch and Dynamic Decoupling (TS+DD) is the combined deployment of TS and DD wherein DD is inserted into slack windows created post gate scheduling, as discussed in Section 4.4.
- Integrated TimeStitch and Dynamic Decoupling w/Heuristic: Integrated TimeStitch and Dynamic Decoupling w/Heuristic (TS+DD(H)) is similar to above, but also incorporates DD insertion according to the heuristic slack threshold.

# 6 EVALUATION

## 6.1 Probability of Success

In Figure 8, we show benefits in terms of POS improvements relative to the ALAP baseline. Benefits shown are in terms of the relative increase in POS. For TS, we show results for TS-SI+C and TS-SI. We also show comparisons to ASAP and Middle. All are detailed in Section 5.4. Applications are ordered by their relative average slack window sizes, and in general, larger slack windows provided greater benefits.

8:20 K. Smith et al.

TS-SI achieves a 50% POS geometric mean improvement, clearly showing the efficiency of the slice+inverse technique in meeting the ideal improvement target. TS-SI+C constrains the slice+inverse technique, so that no SI tuning circuit exceeds the gate depth of the original circuit. Even with this constraint, a mean 38% improvement is obtained, indicating that even under constraint, multiple critical slack windows can be tuned for significant benefits. In comparison, ASAP and Middle achieve no benefit and 15% mean POS improvement, respectively, and observe POS reductions or negligible benefits across many individual benchmarks. While showing some promise, both ASAP and Middle occasionally degrade some benchmarks or provide minimal benefits on others. We can conclude that they are not optimal scheduling solutions. A "one-size-fits-all" approach does not maximize benefits, clear quantitative motivation that specifically-tuned gate positions are preferred. These highlight the benefits of tuning single-qubit gates within slack windows, especially with the practical TS approach that harnesses quantum reversibility in a novel manner for observable quantum circuit gains.

Over the two TS techniques, per-application improvements vary from 5% to 256%. We reason that this variation is caused by the number and sizes of the slack windows, the criticality of the slack window to the circuit, impact of specific gate errors on application fidelity, the input state vectors, as well as general noise characteristics of the machine. Table 1 provides some compiled circuit details. As a note, TS optimizes circuit schedule within slack without increasing the depth or duration of the benchmark.

In Figure 8, plots of CX depth and average slack window size relative to the maximum average slack window are included. It is clear that benefits increase with greater average slack window size. This is important, because slack durations will increase as applications scale and require more SWAPs for qubit communication, as discussed in Section 2.2. We add error bars to the graphs to indicate variation in relative POS benefits from a 1% change to the application's POS. For applications with lower baseline fidelity like the Adder (around 10%), these error bars are longer, but POS benefits are considerable irrespective. This is important as on near-term devices, it is critical in the near to improve the execution of applications with borderline or less than acceptable circuit fidelity. Some variation is expected across runs depending on the particular run's machine characteristics and calibration.

## 6.2 Depth Threshold Sensitivity

In Section 4.3, we motivated the need for restricting the SI tuning circuit depth to the depth of the original circuit. The analysis here involves sweeping thresholds for the depth of a SI circuit from 0 (no slack windows tuned) to 2× the original circuit depth (all slack windows are tunable—equivalent to the unconstrained approach in Section 4.2). Figure 9 shows the POS for the QFT-4 circuit for these different thresholds. The baseline ALAP POS (red line) as well as the depth of the original circuit (blue line) are also shown for comparison. Tuning windows are ordered by tuning circuit size in Figure 9, and TS criteria is satisfied in the first 8 of the 15 windows. Adjusting the target depth threshold of slack tuning can influence the POS. With original circuit depth as the limiting constraint, TS is able to improve the fidelity from 35% to 42%. If machine robustness allowed an unlimited, or at least a 2×, tuning circuit depth, perhaps in the case that the target circuit fell well below coherence bounds, POS jumps to 44%. Depth thresholds can be set based on the machine-application fit. Note that the experimental results suffer from some variation effects of the real machine hence we do not see a strictly monotonically increasing curve.

# 6.3 Comparing Slack Windows

In Figure 10 we show the slack windows of the QEC-5 application, compiled to a five-qubit machine, as a case study. The change in QEC-5 benchmark POS is plotted as gate positions within an isolated

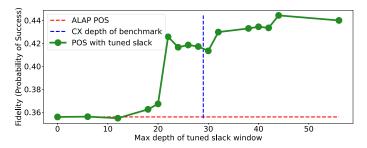


Fig. 9. Threshold sensitivity of window tuning for QFT-4. The red line represents the ALAP POS, the blue line indicates the circuit depth criteria used for TS-SI+C, and the green line describes change in POS as the number of SI tuned slack windows increases.

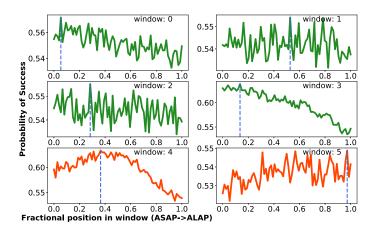


Fig. 10. Slack windows of QEC-5, showing range of POS achieved by tuning each window. Green windows are selected under TS-SI+C while red are rejected. Note that maximum POS for each window differs.

slack window are varied from ASAP (left) to ALAP (right). The windows that are suited to the depth constraint imposed by TS-SI+C are shown in green while the others are in red.

First, it is clear that there are non-negligible POS variations in four of the six windows and all windows have different optimal gate positions. Second, among the green windows, there is considerable benefit in moving to ASAP for window 3. Third, among the red windows, there is considerable benefit for window 4 near the middle. With the TS framework, all local optimums are stitched together to create a final, schedule-optimized circuit. Thus, the benefits of TS-SI+C are considerable over ALAP baseline, and relaxing constraints with TS-SI can produce even greater benefits.

# 6.4 Leveraging Dynamical Decoupling

DD is an established error mitigation technique with similar inspirations as TS. We observe that the mitigation effects of the DD and TS approaches interfere constructively and thus the two can be deployed in a synergistic manner. The benefits can be further improved via intelligent tuning by means of a DD insertion heuristic threshold discussed in Section 4.4.

Figure 11 shows a comparison of TS, DD, DD(H), TS+DD, and TS+DD(H). Note that here, TS is abbreviated for conciseness and corresponds to the constrained TimeStitch approach, TS-SI+C,

8:22 K. Smith et al.

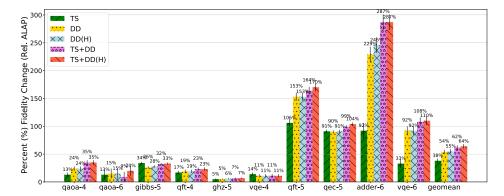


Fig. 11. Comparing POS benefits among TS, DD, and TS+DD variants, relative to ALAP. While DD provides greater benefits than base TS, the error mitigation techniques interfere constructively and the combined approach, TS+DD, performs better than DD or DD(H). DD(H) has the highest average increase in POS. These results were generated with real QC experiments.

as it is practical for real-machine deployment. These are detailed in Section 5.4. All results are normalized to the ALAP baseline.

First, we note that although TS provides significant boosts in benchmark POS, the DD and DD(H) approaches perform equally or better than the constrained TS in all but two benchmarks. The DD approaches are able to achieve geometric mean POS improvements of 54–55% compared to the 38% for constrained TS. The primary reason is that DD can be employed in all slack windows while the constrained TS approach is limited to correcting windows that are allowed by practical circuit depth limitations (Section 4.3). However, constrained TS does not require additional circuit instructions and provides error mitigation with operations already present within the original circuit. It is worth noting from Figure 8 that the unconstrained TS approach achieves a mean 50% POS improvement that is more comparable with the DD benefits and is achieved without adding circuit gates.

While fidelity numbers indicate that DD can outperform TS, it is critical to note that the two approaches are not entirely mitigating the same errors. Additionally, their state refocusing methods differ, coming with their own tradeoffs of either additional tuning or gate overhead. This means that rather than juxtaposing TS and DD as alternatives, there is most opportunity in employing the two techniques in conjunction to maximize the benefits of each. Doing so provides considerable further improvements: TS+DD is able to achieve a 62% POS improvement over ALAP on average, clearly highlighting the synergy between the two techniques. Further, the use of TS and our heuristic threshold that selectively insert DD sequences based on window durations, TS+DD(H), pushes mean improvements to 64% over the ALAP baseline. Overall, the combined approach provides a 19% mean relative improvement over the benefits of DD.

With the combined approach of TS+DD(H), POS improvements range from 7% to a whopping 287%, again dependent on circuit and machine characteristics. These results clearly indicate that combining error mitigation techniques, even beyond those discussed here, can be a considerable fidelity thrust in the NISQ era, resulting in circuits optimized to their full potential.

#### 7 DISCUSSION

The TimeStitch framework targets slack, providing a solution for mitigating decoherence in quantum circuits. In our work, our primary focus is to utilize gates already present in the circuit to implement a technique inspired by Hahn spin-echo. TimeStitch presents the novel contribution of

a slice+inverse tuning mechanism based on quantum reversibility that respects QC frontiers, and single-qubit gates that can be flexibly tuned within slack are targeted for schedule optimization. In addition, our secondary focus is to integrate DD in synergy with the TimeStitch gate-tuning framework. Although TimeStitch requires calibration routines whereas "one-size-fits-all" solutions do not, as described in Section 6, the potential reward of boosted circuit success is well worth the overhead of additional QC runs that can be carefully managed.

TimeStitch focuses on single-qubit gate scheduling within slack. Two-qubit gates could also be eligible for tuning, but they would need careful consideration. To the the best of our knowledge, Hahn spin-echo correction, and DD in general, typically utilizes single-qubit gates. Therefore, tuning the schedule of two-qubit gates would not provide the same type of error mitigation in terms of noise related to qubit state evolution while idling. Two-qubit gate schedule optimization for crosstalk noise mitigation has been seen in the literature [32]. Techniques such as in Reference [32], however, require extending circuit duration through the insertion of barriers. A hallmark of the TimeStitch framework is that it does not extend the duration of the final scheduled circuit. Thus, if two-qubit gates are targeted for schedule optimization, then the only eligible gates would include those not on the circuit's critical path.

#### 7.1 Future Work

Ensuring that quantum optimizations scale along with applications is critical. As discussed in Section 5.3, current TimeStitch slack tuning overheads are manageable as they are contained within a single additional job that must be run before the execution of the final, rescheduled circuit. In near-term QCs where variational noise easily corrupts computation, this small overhead of slice+inverse tuning is trivial, because the average fidelity improvements of +38% on real QCs is significant, outweighing the additional job cost. In some cases, borderline POS values are brought well-above thresholds required for a definite solution because of TimeStitch. We plan to make the TimeStitch framework available so that its optimization techniques can be used by the quantum community at large.

As devices scale, so will applications. As mentioned in Section 2.2, the length of slack in compiled circuits will grow as QC executables increase in width and depth. In this study, circuits were of modest size, enabling thorough slack tuning on all windows. Overheads associated with slack tuning will be kept reasonable by TimeStitch in the future by carefully selecting critical or large windows for tuning. Alternatively, gate position within slack can be sampled using a larger step size. This technique is already implemented in TimeStitch and the user can specify the maximum number of calibration circuits they wish to use. In future work, experiments can be designed to develop heuristics that allow the framework to learn the best windows to target during sampling as well as the ideal search granularity to use. This will likely be influenced by the series of gates that lead up to a window as well as window duration. Additionally, the depth criteria can be strictly enforced to minimize the set of slice+inverse circuits included during optimization. Finally, this work searched for optimum slack schedules though an exhaustive search with a step size dependent on the number of experiments that can fit within a job, but future work will explore refined searching algorithms with lower sampling rates.

A goal of TimeStitch is the optimization of single-qubit gate schedules within circuit slack. When single-qubit gates appear within slack, they are often expressed as a collection of single-qubit operations that are supported by the targeted machine. During single-qubit schedule optimization, TimeStitch currently collapses all operations into one, moving them collectively during SI tuning. While this proves effective for boosting the success of quantum circuits, an interesting avenue for future work would involve independently tuning gates if multiple appeared within circuit slack.

8:24 K. Smith et al.

However, adding additional degrees of freedom to tuning procedures must be done mindfully so that the optimization parameter space does not become prohibitively large to search.

Our investigation focused on IBM QCs but challenges associated with decoherence-related error is not unique to superconducting qubits. Thus, as a final area of future work, we propose the use of TimeStitch to optimally schedule circuits targeted for alternative quantum hardware. Even if other types of qubits are characterized by longer gate and circuit execution times, extending SI tuning procedures, the TS framework would still have potential to provide non-trivial circuit gains.

# 7.2 Related Quantum Proposals

Past work includes the development of methodologies that impact decoherence in quantum circuits by reducing depth and thus overall circuit runtime [11, 28, 29, 46, 55, 60]. These works, however, do not consider variational QC characteristics such as gate error rates and gate durations for their techniques. There also exist frameworks that aim to decrease quantum circuit noise by taking device calibration data into consideration to improve program success [31, 49, 57], but these techniques do not implement optimizations that take advantage of slack time in circuits. Next, optimizing schedulers exist that mitigate noise associated with crosstalk by considering device properties [14, 32]. In addition, the methods in Reference [58] take advantage of quantum circuit slack but focus on the qubit mapping problem rather than error reduction on real QCs. Further, the benefits of our method complement other indirect decoherence mitigation approaches [25, 31]. Finally, related work exists for both quantum reversibility and DD applied in quantum circuit optimization. These proposals and their relation to the TimeStitch framework are discussed in detail in Sections 2.3 and 3.3, respectively.

# 7.3 Exploiting Slack in Classical Computing

Optimizing circuit slack is not limited to quantum computation. With the critical need for maximizing performance under strict power budgets in the classical world, exploiting any form of idle time or slack in classical computing has been a popular theme of research over past decades. At the circuit level, slack in a clock cycle can occur in the presence of conservative timing guardbands. These have been exploited with multiple better-than-worse-case approaches [16, 22, 40–42, 51]. Similarly, at the micro-architecture level, periods of time with less or no-activity can help save power at no additional performance costs. These are often exploited via power/clock gating, multi threading [52], instruction rescheduling [18], and so on.

## 8 CONCLUSION

Reducing the impact of decoherence is critical for substantial advancements on near-term QCs. The unintentional coupling of qubits to their environment, and each other, adds significant noise to computation, and improved methods to combat decoherence are required to boost the performance of quantum algorithms on real machines. This article presents a novel technique that takes advantage of a largely unexplored space of quantum circuit slack, opening up a new domain of exploration.

Quantum circuit slack will only become more prevalent in time. Here, slack tuning improves the fidelity of compiled quantum circuits without either increasing total gate count or introducing circuit partitioning that increases circuit duration. By exploiting quantum reversibility and by constraining tuning circuits to the depth of the original application, we propose a practical design suited for a variety of applications and quantum machines, especially applications of low fidelity that are critical to improve. We evaluated our proposal TimeStitch on real quantum machines and on benchmarks that are critical to real-world quantum usecases. We additionally offer insights on challenges and optimizations suited to realistic deployment.

#### REFERENCES

- [1] IBM. IBM Quantum Experience. Retrieved November 18, 2020 from https://quantum-computing.ibm.com.
- [2] IBM. IBM Quantum SU2 ansatz. Retrieved November 18, 2020 from https://qiskit.org/documentation/stubs/qiskit.circuit.library.EfficientSU2.html.
- [3] IBM. IBM Quantum Systems. Retrieved from November 18, 2020 from https://quantum-computing.ibm.com/services?systems=all.
- [4] Héctor Abraham, Adu Offei, Rochisha Agarwal, Ismail Yunus Akhalwaya, Gadi Aleksandrowicz, Thomas Alexander, Matthew Amy, et al. 2019. Qiskit: An Open-source Framework for Quantum Computing. https://github.com/Qiskit/qiskit/blob/master/Qiskit.bib.
- [5] Frank Arute, Kunal Arya, Ryan Babbush, Dave Bacon, Joseph C. Bardin, Rami Barends, Andreas Bengtsson, Sergio Boixo, Michael Broughton, Bob B. Buckley, et al. 2020. Observation of separated dynamics of charge and spin in the fermi-hubbard model. arXiv:2010.07965. Retrieved from https://arxiv.org/abs/2010.07965.
- [6] Jacob Biamonte, Peter Wittek, Nicola Pancotti, Patrick Rebentrost, Nathan Wiebe, and Seth Lloyd. 2017. Quantum machine learning. Nature 549, 7671 (2017), 195–202.
- [7] M. J. Biercuk, A. C. Doherty, and H. Uys. 2011. Dynamical decoupling sequence construction as a filter-design problem. *J. Phys. B* 44, 15 (2011), 154002.
- [8] A. Robert Calderbank, Eric M. Rains, P. M. Shor, and Neil J. A. Sloane. 1998. Quantum error correction via codes over GF (4). *IEEE Trans. Inf. Theory* 44, 4 (1998), 1369–1387.
- [9] Christopher Chamberland, Guanyu Zhu, Theodore J. Yoder, Jared B. Hertzberg, and Andrew W. Cross. 2020. Topological and subsystem codes on low-degree graphs with flag qubits. *Phys. Rev. X* 10, 1 (2020), 011022.
- [10] Zijun Chen, Kevin J. Satzinger, Juan Atalaya, Alexander N. Korotkov, Andrew Dunsworth, Daniel Sank, Chris Quintana, Matt McEwen, Rami Barends, Paul V. Klimov, et al. 2021. Exponential Suppression of Bit or Phase Flip Errors with Repetitive Error Correction. *Nature* 595, 7867 (2021), 383–387.
- [11] Andrew M. Childs, Eddie Schoute, and Cem M. Unsal. 2019. Circuit Transformations for Quantum Architectures. In Proceedings of the 14th Conference on the Theory of Quantum Computation, Communication and Cryptography (TQC'19), Wim van Dam and Laura Mancinska (Eds.). Vol. 135, Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 3:1–3:24. DOI:10.4230/LIPIcs.TQC.2019.3
- [12] Steven A. Cuccaro, Thomas G. Draper, Samuel A. Kutin, and David Petrie Moulton. 2004. A new quantum ripple-carry addition circuit. arXiv preprint quant-ph/0410184.
- [13] Poulami Das, Swamit Tannu, Siddharth Dangwal, and Moinuddin Qureshi. 2021. ADAPT: Mitigating Idling Errors in Qubits via Ddaptive Dynamical Decoupling. In *Proceedings of the 54th Annual IEEE/ACM International Symposium on Microarchitecture*. 950–962.
- [14] Yongshan Ding, Pranav Gokhale, Sophia Fuhui Lin, Richard Rines, Thomas Propson, and Frederic T. Chong. 2020. Systematic Crosstalk Mitigation for Superconducting Qubits via Frequency-aware Compilation. In Proceedings 53rd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO'20). IEEE, 201–204.
- [15] David P. DiVincenzo and Peter W. Shor. 1996. Fault-tolerant error correction with efficient quantum codes. *Phys. Rev. Lett.* 77, 15 (1996), 3260.
- [16] D. Ernst, Nam Sung Kim, S. Das, S. Pant, R. Rao, Toan Pham, C. Ziesler, D. Blaauw, T. Austin, K. Flautner, and T. Mudge. 2003. Razor: A low-power pipeline based on circuit-level timing speculation. In *Proceedings of the 36th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO'36)*. 7–18. https://doi.org/10.1109/MICRO.2003. 1253179
- [17] Edward Farhi, Jeffrey Goldstone, and Sam Gutmann. 2014. A Quantum Approximate Optimization Algorithm. arXiv:1411.4028. Retrieved from https://arxiv.org/abs/1411.4028.
- [18] Brian Fields, Rastislav Bodík, and Mark D. Hill. 2002. Slack: Maximizing performance under technological constraints. In Proceedings of the 29th Annual International Symposium on Computer Architecture (ISCA'02). IEEE Computer Society, 47–58.
- [19] Austin G. Fowler, Matteo Mariantoni, John M. Martinis, and Andrew N. Cleland. 2012. Surface codes: Towards practical large-scale quantum computation. *Phys. Rev. A* 86, 3 (2012), 032324.
- [20] Tudor Giurgica-Tiron, Yousef Hindy, Ryan LaRose, Andrea Mari, and William J. Zeng. 2020. Digital zero noise extrapolation for quantum error mitigation. In Proceedings of the IEEE International Conference on Quantum Computing and Engineering (QCE'20). IEEE, 306–316.
- [21] Daniel M. Greenberger, Michael A. Horne, and Anton Zeilinger. 1989. Going beyond Bell's theorem. In *Bell's Theorem*, *Quantum Theory and Conceptions of the Universe*. Springer, 69–72.
- [22] Meeta Gupta, Jude A. Rivers, Pradip Bose, Gu-Yeon Wei, and David Brooks. 2009. Tribeca: Design for PVT variations with local recovery and fine-grained adaptation. In *Proceedings of the 42nd Annual IEEE/ACM International Symposium on Microarchitecture*. 435–446.

8:26 K. Smith et al.

[23] Petar Jurcevic, Ali Javadi-Abhari, Lev S. Bishop, Isaac Lauer, Daniela Borgorin, Markus Brink, Lauren Capelluto, Oktay Gunluk, Toshinari Itoko, Naoki Kanazawa, et al. 2021. Demonstration of quantum volume 64 on a superconducting quantum computing system. Quant. Sci. Technol. 6, 2 (2021), 025020.

- [24] Abhinav Kandala, Antonio Mezzacapo, Kristan Temme, Maika Takita, Markus Brink, Jerry M. Chow, and Jay M. Gambetta. 2017. Hardware-efficient variational quantum eigensolver for small molecules and quantum magnets. Nature 549, 7671 (2017), 242–246.
- [25] Gushu Li, Yufei Ding, and Yuan Xie. 2019. Tackling the qubit mapping problem for NISQ-era quantum devices. In Proceedings of the 24th International Conference on Architectural Support for Programming Languages and Operating Systems. 1001–1014.
- [26] Gang-Qin Liu, Hoi Chun Po, Jiangfeng Du, Ren-Bao Liu, and Xin-Yu Pan. 2013. Noise-resilient quantum evolution steered by dynamical decoupling. *Nat. Commun.* 4, 1 (August 2013), 1–9. https://doi.org/10.1038/ncomms3254
- [27] Ryan Mandelbaum. 2021. Five Years Ago Today, We Put the First Quantum Computer on the Cloud. Here's How We Did It. Retrieved October 19, 2021 from https://research.ibm.com/blog/quantum-five-years.
- [28] Dmitri Maslov, Sean M. Falconer, and Michele Mosca. 2008. Quantum circuit placement. IEEE Trans. Comput.-Aid. Des. Integr. Circ. Syst. 27, 4 (2008), 752–763.
- [29] Tzvetan S. Metodi, Darshan D. Thaker, Andrew W. Cross, Frederic T. Chong, and Isaac L. Chuang. 2006. Scheduling physical operations in a quantum information processor. In *Quantum Information and Computation IV*, Vol. 6244. International Society for Optics and Photonics, 62440T.
- [30] Nikolaj Moll, Panagiotis Barkoutsos, Lev S. Bishop, Jerry M. Chow, Andrew Cross, Daniel J. Egger, Stefan Filipp, Andreas Fuhrer, Jay M. Gambetta, Marc Ganzhorn, et al. 2018. Quantum optimization using variational algorithms on near-term quantum devices. *Quant. Sci. Technol.* 3, 3 (2018), 030503.
- [31] Prakash Murali, Jonathan M. Baker, Ali Javadi-Abhari, Frederic T. Chong, and Margaret Martonosi. 2019. Noise-adaptive compiler mappings for noisy intermediate-scale quantum computers. In Proceedings of the 24th International Conference on Architectural Support for Programming Languages and Operating Systems. 1015–1029.
- [32] Prakash Murali, David C. McKay, Margaret Martonosi, and Ali Javadi-Abhari. 2020. Software mitigation of crosstalk on noisy intermediate-scale quantum computers. In Proceedings of the 25th International Conference on Architectural Support for Programming Languages and Operating Systems. 1001–1016.
- [33] Paul Nation, Hanhee Paik, Andrew Cross, and Zaira Nazario. 2021. The IBM Quantum Heavy Hex Lattice. Retrieved August 11, 2021 from https://www.research.ibm.com/blog/heavy-hex-lattice.
- [34] Michael A. Nielsen and Isaac Chuang. 2010. Quantum Computation and Quantum Information. Cambridge University Press.
- [35] Tirthak Patel and Devesh Tiwari. 2021. Qraft: Reverse your Quantum circuit and know the correct program output. Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems. 443–455.
- [36] Alberto Peruzzo, Jarrod McClean, Peter Shadbolt, Man-Hong Yung, Xiao-Qi Zhou, Peter J. Love, Alán Aspuru-Guzik, and Jeremy L. O'brien. 2014. A variational eigenvalue solver on a photonic quantum processor. *Nat. Commun.* 5, 1 (2014), 1–7.
- [37] Bibek Pokharel, Namit Anand, Benjamin Fortman, and Daniel A. Lidar. 2018. Demonstration of fidelity improvement using dynamical decoupling with superconducting qubits. Phys. Rev. Lett. 121, 22 (2018), 220502.
- [38] John Preskill. 2018. Quantum computing in the NISQ era and beyond. Quantum 2 (2018), 79.
- [39] Gregory Quiroz and Daniel A. Lidar. 2011. Quadratic dynamical decoupling with nonuniform error suppression. Phys. Rev. A 84, 4 (2011), 042328.
- [40] Gokul Subramanian Ravi. 2020. Integrating Computing Systems from the Gates Up: Breaking the Clock Abstraction. The University of Wisconsin—Madison.
- [41] Gokul Subramanian Ravi and M. Lipasti. 2018. Aggressive slack recycling via transparent pipelines. In *Proceedings of the International Symposium on Low Power Electronics and Design*. 1–6.
- [42] G. S. Ravi and M. Lipasti. 2019. Recycling data slack in out-of-order cores. In Proceedings of the IEEE International Symposium on High Performance Computer Architecture (HPCA'19). 545–557. https://doi.org/10.1109/HPCA.2019.00065
- [43] Gokul Subramanian Ravi, Kaitlin N. Smith, Prakash Murali, and Frederic T. Chong. 2021. Adaptive job and resource management for the growing quantum cloud. In Proceedings of the IEEE International Conference on Quantum Computing and Engineering. IEEE, 301–312.
- [44] Joschka Roffe. 2019. Quantum error correction: An introductory guide. Contemp. Phys. 60, 3 (July 2019), 226–245. https://doi.org/10.1080/00107514.2019.1667078
- [45] Peter W. Shor. 1999. Polynomial-time algorithms for prime factorization and discrete logarithms on a quantum computer. SIAM Rev. 41, 2 (1999), 303–332.
- [46] Marcos Yukio Siraichi, Vinícius Fernandes dos Santos, Sylvain Collange, and Fernando Magno Quintão Pereira. 2018. Qubit allocation. In *Proceedings of the International Symposium on Code Generation and Optimization*. 113–125.

- [47] Alexandre M. Souza, Gonzalo A. Álvarez, and Dieter Suter. 2012. Robust dynamical decoupling. *Philos. Trans. Roy. Soc. A: Math. Phys. Eng. Sci.* 370, 1976 (2012), 4748–4769.
- [48] Alexandre M. Souza, Gonzalo A. Álvarez, and Dieter Suter. 2012. Robust dynamical decoupling. *Philos. Trans. Roy. Soc. A: Math. Phys. Eng. Sci.* 370, 1976 (October 2012), 4748–4769. https://doi.org/10.1098/rsta.2011.0355
- [49] Swamit S. Tannu and Moinuddin K. Qureshi. 2019. Not All qubits are created equal: A case for variability-aware policies for NISQ-era quantum computers. In Proceedings of the 24th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'19). Association for Computing Machinery, New York, NY, 987–999. https://doi.org/10.1145/3297858.3304007
- [50] Kristan Temme, Sergey Bravyi, and Jay M. Gambetta. 2017. Error mitigation for short-depth quantum circuits. Phys. Rev. Lett. 119, 18 (2017), 180509.
- [51] Abhishek Tiwari, Smruti R. Sarangi, and Josep Torrellas. 2007. ReCycle: Pipeline adaptation to tolerate process variation. In *Proceedings of the International Symposium on Computer Architecture (ISCA'07)*. 323–334.
- [52] Dean M. Tullsen, Susan J. Eggers, and Henry M. Levy. 1995. Simultaneous multithreading: Maximizing on-chip parallelism. In Proceedings of the 22nd Annual International Symposium on Computer Architecture (ISCA'95). Association for Computing Machinery, New York, NY, 392–403. https://doi.org/10.1145/223982.224449
- [53] A. M. Tyryshkin, Zhi-Hui Wang, Wenxian Zhang, E. E. Haller, J. W. Ager, V. V. Dobrovitski, and S. A. Lyon. 2010. Dynamical decoupling in the presence of realistic pulse errors. arXiv:1011.1903. Retrieved from https://arxiv.org/abs/1011.1903.
- [54] Götz S. Uhrig. 2007. Keeping a quantum bit alive by optimized  $\pi$ -pulse sequences. *Phys. Rev. Lett.* 98, 10 (2007), 100504.
- [55] Davide Venturelli, Minh Do, Eleanor Rieffel, and Jeremy Frank. 2018. Compiling quantum circuits to realistic hardware architectures using temporal planners. Quant. Sci. Technol. 3, 2 (2018), 025004.
- [56] Lorenza Viola, Emanuel Knill, and Seth Lloyd. 1999. Dynamical decoupling of open quantum systems. Phys. Rev. Lett. 82, 12 (1999), 2417.
- [57] Christophe Vuillot. 2018. Is error detection helpful on IBM 5Q chips? *Quantum Information and Computation* 18, 11–12 (2018), 949–964.
- [58] Chi Zhang, Yanhao Chen, Yuwei Jin, Wonsun Ahn, Youtao Zhang, and Eddy Z. Zhang. 2020. SlackQ: Approaching the qubit mapping problem with a slack-aware swap insertion scheme. arXiv:2009.02346. Retrieved from https://arxiv. org/abs/2009.02346.
- [59] Christa Zoufal, Aurélien Lucchi, and Stefan Woerner. 2021. Variational quantum Boltzmann machines. Quant. Mach. Intell. 3, 1 (February 2021), 1–15. https://doi.org/10.1007/s42484-020-00033-7
- [60] Alwin Zulehner and Robert Wille. 2019. Compiling SU (4) quantum circuits to IBM QX architectures. In Proceedings of the 24th Asia and South Pacific Design Automation Conference. 185–190.

Received 26 October 2021; revised 21 May 2022; accepted 22 June 2022